



Peer to peer size estimation in large and dynamic networks: A comparative study

Erwan Le Merrer, Anne-Marie Kermarrec, Laurent Massoulié

► To cite this version:

Erwan Le Merrer, Anne-Marie Kermarrec, Laurent Massoulié. Peer to peer size estimation in large and dynamic networks: A comparative study. HPDC-15, Jun 2006, Paris. inria-00080652v3

HAL Id: inria-00080652

<https://inria.hal.science/inria-00080652v3>

Submitted on 21 Jul 2006

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Peer to peer size estimation in large and dynamic networks: A comparative study

Erwan Le Merrer [†], Anne-Marie Kermarrec [‡] & Laurent Massoulié [★]

[†] France Telecom R&D, Lannion, France and IRISA

[‡] INRIA/IRISA, Rennes, France

[★] Microsoft Research, Cambridge, UK

E-mail: erwan.lemerrer@orange-ft.com, Anne-Marie.Kermarrec@irisa.fr, lmassoul@microsoft.com

Abstract

As the size of distributed systems keeps growing, the peer to peer communication paradigm has been identified as the key to scalability. Peer to peer overlay networks are characterized by their self-organizing capabilities, resilience to failure and fully decentralized control. In a peer to peer overlay, no entity has a global knowledge of the system. As much as this property is essential to ensure the scalability, monitoring the system under such circumstances is a complex task. Yet, estimating the size of the system is a core functionality for many distributed applications to parameter setting or monitoring purposes. In this paper, we propose a comparative study between three algorithms that estimate in a fully decentralized way the size of a peer to peer overlay. Candidate approaches are generally applicable irrespective of the underlying structure of the peer to peer overlay. The paper reports the head to head comparison of estimation system size algorithms. The simulations have been conducted using the same simulation framework and inputs and highlight the differences in cost and accuracy of the estimation between the algorithms both in static and dynamic settings.

I. Introduction

The past decade has been marked by a tremendous shift in the scale of distributed systems. One of the most striking consequences is that most of traditional algorithms are no longer relevant in such a context. To cope with an increasing number of participants, the peer to peer communication paradigm has been imposed as a key to fill this gap and many algorithms relies on such a model both in industry and academia.

Peer to peer systems are scalable, self-organizing and

resilient to failures. They provide an attractive support for many applications from file sharing systems (e.g. Kazaa, Gnutella) to video over IP applications (e.g. CoolStreaming), or voice on IP applications (e.g. Skype). A peer to peer overlay network connects peers in a logical manner on top of IP so that each peer is aware of a small subset of the network. Yet, the sum of individual decisions based on a restricted knowledge makes the system converging towards global properties.

As much as the fully decentralized nature of a peer to peer system is essential to the scalability, it makes the system monitoring much more complex. Obtaining global statistics on the system becomes a hard issue. However, such statistics might be crucial to adjust the parameters of the system as the structure and characteristics change over time. System size estimation is one of the most crucial statistic to be collected on a system since it can be used for many purposes such as parameter value setting or network monitoring. For example, some peer to peer protocols require this estimation as an input parameter: the constant degree of the Viceroy network [12] requires this information to choose a *level* for an incoming peer. Some gossip-based protocols use this information to compute the number of gossip targets of a message to reach given propagation probabilities [7],[4].

Whereas the number of participants can easily be tracked by a set of servers in server-based systems (e.g. Naspter), it is no longer obvious in a fully unstructured network where each node is connected to a set of random neighbors (e.g. Gnutella, Freenet). In the context of structured overlay networks, several algorithms relying on the actual structure of the network have been proposed ([17], [11], [13], [14]). Such algorithms exploit the fact that node identifiers are uniformly assigned at random. The size estimation may then be directly inferred from the observation of the density of identifiers that fall into a given subset of the global identifier space.

While those methods provide good approximation of the system size, their applicability is strictly limited to those identifier-based overlay networks. Instead, in this paper, we consider generally applicable solutions irrespective of the underlying network topology. The main motivation of this paper is to compare several generally applicable algorithms to estimate the system size in both static and dynamic contexts in an attempt to help application developers to choose the best strategy for a given setting/cost/accuracy. To this end, we have carefully chosen three candidate algorithms, each representative of a class of approaches, for the head to head comparison. The first algorithm relies on random walks, using the inverted birthday paradox, the second candidate is based on probabilistic polling while the last one relies on epidemic-based information aggregation. To the best of our knowledge such a comparison is the first of its kind. The only paper we are aware of comparing different fully distributed algorithms is [17] comparing the Hops Sampling method, discussed in this paper, with an interval density based approach.

We implemented the three candidates algorithms and simulated their behavior against the same inputs for a fair comparison. In this paper, we report the simulation results and compare the methods along their accuracy, overhead and reactivity to changes in the environment.

The rest of the paper is structured as follow. In the next section, we provide some background of generic solution to estimate in a decentralized way the size of a large system. In Section III, we successively describe the three candidate algorithms. These methods are then evaluated and compared and the results are reported in Section IV. We discuss the pros and cons of the compared candidates based on the simulation results in Section V and conclude in Section VI.

II. Background on generic counting algorithms

In this section, we provide some background on the three main generic distributed counting approaches we have identified.

The algorithms belonging to the first class rely on called *probabilistic polling* or *randomized report* techniques. In such approach, the basic idea is to probe the network in a probabilistic way and to infer the size of the system based on the set of replies. To this end, the initiator node broadcasts a message to all nodes in the overlay, and waits for responses. Upon receipt of a broadcast message, the nodes send back a response with a probability depending on the probability parameter set in the broadcast message [2], [6], or based on their distance to the initiator node [11] for example. In this paper, we chose the last one as a candidate for comparison; the fact of using node's distance to the

initiator node could lower message overhead compared to simple probabilistic response, as fewer "far nodes" should reply with messages that will cross an important part of the overlay. This algorithm has also been evaluated in [11] and will permit us by results comparison to validate our simulator approach.

The second class of approaches relies on epidemic algorithms to permanently gather information on the system by constantly exchanging information between peers, mostly in a random manner [9], [8], [19]. Epidemic algorithms have recently received a lot of attention due to their great resilience to failure; they also have been used for aggregation applications [18]. The algorithm presented in [9] is based on a *push-pull anti-entropy epidemic* protocol. In this protocol, each peer periodically selects another peer at random to exchange information with. During this epidemic propagation, an *aggregation* of values will gradually take place, hosted by each cooperating overlay node, thus leading to a global value aggregate representing the targeted system size. This method has proven to be accurate and to converge pretty quickly. We chose this algorithm as a candidate.

Finally, the last class gathers the approaches relying on random walks. For example, [2] proposes a random *increasing* walk on the topology, in an overlay where each node owns an identifier. Starting from the node with the smallest *id*, the message is forwarded toward higher *id* nodes. This method however suffers from the same drawbacks as density-based approaches since it relies on identifier assignment. [15] first introduces the *Random Tour* method, based on an emulation of the return time of a random walk to the initiating node. The second method proposes [15], the *Sample&Collide* algorithm, uses random walks to produce an unbiased sampling (uniformly at random) of overlay nodes. Those samples are then used to compute the system size according to the *inverted birthday paradox* introduced in [2]. Given the results provided in [15], the overhead of the *Sample&Collide* algorithm is much lower than the one of Random Tour. Therefore, we chose the *Sample&Collide* algorithm as the candidate of this class.

III. Candidate algorithms overviews

A. Sample&Collide

The *Sample&Collide* method [15] builds upon the *inverted birthday paradox* method introduced in [2]. The birthday paradox states that the probability of having two people in a room with the same birthday is at least 1/2 for a group of 23 people or more. More generally, detailed evaluations of the probability $p(N, K)$ of having at least two persons with the same birthday, drawn from a set of

$N = 365$ dates (ignoring leap years for simplicity), within a group of K people, are available.

The inverted birthday paradox method relies on inverting such evaluations, to determine the probability distribution of the random number $X(N)$ of people one needs to have, adding them one at a time, until two of them share the same birthday. Clearly, the probability that $X(N)$ equals K is given by $p(N, K) - p(N, K - 1)$. It turns out that for large N , $X(N)$ is concentrated around $\sqrt{2N}$.

The inverted birthday paradox method relies on this fact, and proceeds as follows. The dates are the peers in the system. Random samples of peers are obtained consecutively (these correspond to the people) until two samples coincide. The number of samples X drawn till this happens is then used to provide the estimate \hat{N} of the number of peers, $\hat{N} = X^2/2$.

This algorithm heavily relies on the correctness of the sampling method used. The *Sample&Collide* method improves upon the *inverted birthday paradox* method in two respects. First, the sampling technique used in [15] is asymptotically unbiased, on arbitrary graphs, while previously proposed sampling schemes are typically biased in graphs with heterogeneous node degrees. Second, the samples are used more efficiently in *Sample&Collide* than in inverted birthday paradox, thereby yielding a better accuracy/cost trade-off.

The uniform peer sampling technique which produces unbiased samples is implemented as follows: the initiator node sets a predefined value $T > 0$. This value is then sent to a neighbor chosen uniformly at random. Each node receiving the message first picks a random number U , uniformly distributed on $[0, 1]$; it then simply decrements T by $-\log(U)/d_i$ (d_i is the degree of the current node), and forwards the message to a neighbor, if $T > 0$. Otherwise the current node is the sample node, and it returns its *id* to the initiator.

Asymptotic unbiasedness means that the distribution of the returned sample approaches the uniform distribution as the parameter T increases. The expansion properties of the graph influence how large T should be selected in order to have negligible bias.

The accuracy/overhead tradeoff of this algorithm relies on the control parameter l which determines the number of newly sampled nodes that have already been observed. The greater l , the more precise is the estimation, but the higher is the generated overhead.

The details of the algorithm, the analysis and simulation results are available in [15] as well as a comparison with the *Random tour* approach. Results show that the *Sample&Collide* approach improves upon the *Random Tour* and justify the choice of the former in this comparative study.

B. HopsSampling

The second candidate approach we chose is called the *HopsSampling* and represents in this comparative study the class of probabilistic polling approaches. Two versions of the Hops Sampling algorithm are presented in [17] and [11]: the *gossipSample* and *minHopsReporting* heuristics.

As a first step to choose the candidate approach for the comparison, we tried to reproduce the results obtained in [17] and [11]. We were able to reproduce the results obtained using the *minHopsReporting* heuristic. The results we obtained by implementing the pseudo-code in [17] (*gossipSample* heuristic) somehow led to less accurate results. Therefore for the sake of fairness, and after discussions with the authors [16], we decided to choose the system size estimation algorithm relying on the *minHopsReporting* heuristic.

This algorithm is based on the probabilistic polling technique where an initiator spreads messages in the network and estimates the system size based on the replies it gets back. These messages are sent back in a probabilistic manner. The HopsSampling algorithm, evaluated in this paper works as follows: an initiator node spreads by gossip a message across the network. Initially this message contains *hopCount* set to 0. At each traversed node, this value is incremented; the lowest *hopCount* value received by a node is remembered and represents its distance from the count initiator. Then, depending on each node's distance (in order to avoid massive flood towards the initiator), there is a probabilistic response message sent to the initiator node: (i) if *hopCount* < *minHopsReporting*, a response is set with probability 1, else (ii) the response is sent with probability $\frac{1}{\frac{gossipTo}{hopCount} - minHopsReporting}$. For each message count received from nodes at a certain distance, the initiator needs to multiply it by the percentage of peers in the network they represents. For example, if *minHopsReporting* = 2, only 25% of nodes with distance 4 will report back. The number of collected responses will thus be multiplied by $\frac{1}{0.25}$, i.e. 4. The system size is then extrapolated from these collected returns.

This algorithm, and particularly the initial gossip spread, works with a priori fixed parameter values discussed in [17], [11]. We have also talked about those parameters value in [16] to ensure a fair comparison.

C. Gossip-based Aggregation

The last chosen approach for comparison is representative of the class of epidemic-based approaches. Epidemic-based protocols have recently received a lot of attention given their scalability properties. In an epidemic protocol, each peer periodically exchanges information with one of its neighbor picked at random. In [9], a gossip-based

aggregation protocol is presented and the results, that we were able to reproduce, show that this approach provides accurate results.

This approach [9] relies on the following statement: if exactly one node of the system holds a value equal to 1, and all the other values are equal to 0, the average is $1/N$. The system size could thus be directly computed. To run this algorithm, an initiator should take the value equal to 1, and start gossiping; the reached nodes participate to the process by setting their value to 0. At each predefined cycle, each node in the network chooses one of its neighbor at random and swaps its estimation parameter. The contacted node does the same (push/pull heuristic of [9]). Both nodes then recompute their estimation as follows:

$$Estimation \leftarrow \frac{Estimation + Neighbor's_Estimation}{2}$$

To provide correct estimations, this algorithm needs to wait a certain number of *rounds* to elapse before computing the size estimation; this period is the required time for the gossip to propagate in the whole overlay and for the values to converge.

This method converges toward exact system size in a stable system. More details may be found in [9].

D. Summary

We have chosen three candidate algorithms, that we believe are representative of three classes of generic solutions to estimate the size of a large-scale peer to peer system in a fully decentralized way. They are applicable to any peer to peer overlay in which each node is connected to a set of *random* neighbors. Most of structured overlays, such as e.g. Pastry, fall into this category.

The goal of the comparative study is then to run simulations in the same framework in order to compare those methods not only in terms of accuracy but overhead and resilience to dynamic changes.

IV. Evaluation

A. Experimental Setup

One of the main metrics along which we compare the different algorithms is their scalability in terms of overhead, as well as their accuracy as the system size varies. Given the considered scale, which could not realistically be reached in a lab testing environment, we evaluated them using a discrete event simulator, able to simulate static and dynamic network configurations. The simulator counts the messages over the network. It does not model the physical network topology nor the queuing delays and packet losses.

The first step of the simulation, as mentioned in the previous section, was to reproduce the results obtained by the authors of each algorithm. This step has been successfully reached for the three considered algorithms and we now compare them in the same simulation environment.

The simulations were run on unstructured peer to peer networks. To this end, we built the peer to peer overlay as a random graph, where each node was provided with a set of k neighbors, chosen uniformly at random among the whole set of peers. We do not consider in this paper the actual construction of such graphs but several approaches exist to build such peer to peer overlay in practice [10]. In order to be as close as possible from real world implementations, the graphs we used in the evaluation were not fully homogeneous graphs where each node had the same number of neighbors. Instead, each node owns a random number of neighbors, below a given threshold. This reflects for instance heterogeneity between nodes in terms of number of incoming and outgoing links capacities. We also run some tests in the context of homogeneous graphs. This parameter consistently improved all algorithms. Therefore, we chose the worst case setting and present the results of experiments conducted on heterogeneous networks.

Graphs construction: In test configurations, each node has a number of neighbors varying between 1 and a fixed max value. At the beginning of the construction process, all nodes are present in the overlay (thus available to be chosen as neighbors). Nodes are taken one by one to be wired: the current node first chooses uniformly at random its current number of neighbors, and fills its view with again uniformly at random selected nodes as a neighbors, that do not already have the max fixed value (otherwise other random nodes are chosen). We used 10 neighbors max for the simulations, which leads in both overlay sizes to an average of approximately 7.2. This value, over $\log(N)$ (N being the size of the system, and \log is \log_{10}), ensures graph connectivity [5]: *in order for a network to stay connected with constant probability when all nodes fail with probability 1/2, some nodes must have degree $O(\log n)$; and if no node is substantially above-average degree, the lemma can be strengthened: the average degree must be $O(\log n)$.*

We also consider a *scale free* graph to test the algorithms on. It is based on the Barabási and Albert model [1] with growth and preferential attachment; the experiments could be found at the end of Section IV-C.

The links between nodes are assumed to be bidirectional (whenever a node contacts another one, the reached node also has knowledge of communication initiator's existence and keeps a link back to the contact node), leading to an undirected graph.

We scaled up to 1,000,000 node graphs, but dynamic environment was created on 100,000 node graphs for

practical considerations. Node removals introduce a loss of connectivity for the remaining nodes; in this simulation, the nodes that have lost one or several neighbors do not create new links with other nodes.

The rest of the section is organized as follows: We first describe the evaluation criteria. In Section IV-C, we report the results obtained in static contexts and in Section IV-D, we introduce some dynamics in the system to evaluate the various algorithms capabilities to provide reactive estimations. Section IV-E provides some overhead results.

B. Evaluation criteria

We used several metrics to evaluate the different approaches to estimate in a distributed fashion the size of a large-scale dynamic system.

a) Accuracy: One of the most important criterion is the quality of the estimate and the speed to which this accuracy is reached. This is an important parameter for the use of the system size information. For some applications, a quick approximation might be more appropriate than a more accurate one which would take much longer to compute, or be more expensive to obtain.

b) Reactivity to changes: A second characteristic we want to compare the approaches against is their ability to react to changes in the environment. To this end, we compute the time to react to a growth or increase of the number of peers in the system.

c) Overhead: The third important criterion is the overhead of each approach. In this paper, we compute the overhead as the number of messages required to compute the system size. This information could be extrapolated by forthcoming designers to approximate the bandwidth and the computational resources that will be needed by the application.

C. Simulation results in static contexts

In the following set of experiments, the algorithms run on an already constructed overlay where peers are connected to a random graph, which size remains static throughout the experiments. On the following figures, the system size is normalized to 100 to enable us to express the quality of the estimation in terms of percentage.

d) Sample&Collide: The quality of *Sample&Collide*'s estimations depends upon the value of the T and l parameters. In the rest of the experiments we set T to 10; this value is sufficient for an accurate sampling [15]. Figures 1 and 2 depict the quality (on the y axis) of an estimation respectively in a 100,000 and 1,000,000 node configuration depending on the number of estimations computed over time (x axis). The dashed

curve represents a *oneShot* estimation which is a single non averaged estimation while the continuous curve, *last10runs*, is the average of the 10 last estimations. In these experiments we set l to 200, we will consider other values later in the evaluation.

Our first observation is that the results show that the *oneShot* curve remains most of the time in a 10% precision window, with some peaks between 10 and 20%. The *last10runs* one remains within 3 or 4% of the exact value. Obviously *oneShot* has a very low overhead as compared to the *last10run* one. Therefore, in the rest of the paper, the *Sample&Collide* approach will be related to the *oneShot* configuration.

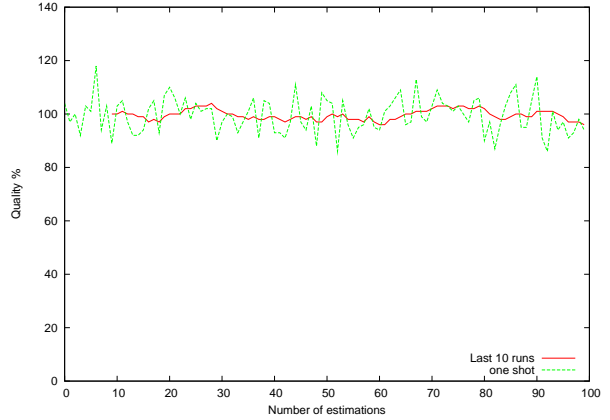


Fig. 1. *Sample&Collide*: oneShot and last10runs heuristic with $l=200$, 100,000 node network, static environment

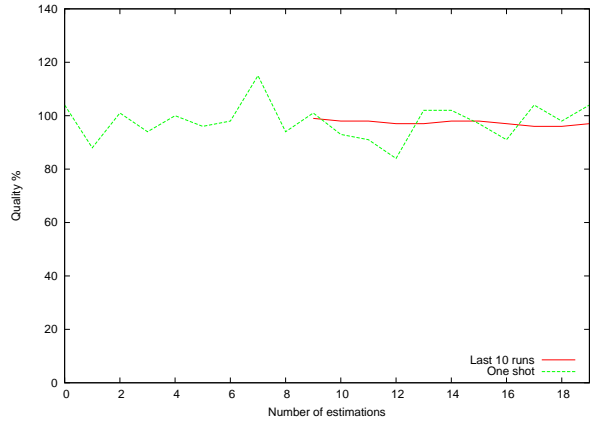


Fig. 2. *Sample&Collide*: oneShot and last10runs heuristic with $l=200$, 1,000,000 node network

e) HopsSampling: We evaluated this algorithm first in small size networks in order to reproduce the results presented in [11], and observed similar results. We used in our implementation the gossip algorithm presented in the original papers, with the following values for the parameters: $gossipTo = 2$, $gossipFor = 1$, $gossipUntil = 1$, $minHopsReporting = 5$ ([17], [16]).

Figures 3 and 4 depict the quality of the estimation (y axis), depending on the number of runs (x axis), respec-

tively in a 100,000 and 1,000,000 node configuration. As in *Sample&Collide*, we report the *oneShot* and average over the *last10runs* results. We observe that the algorithm scales well. We also observe not surprisingly that the *last10runs* heuristic provides a less noisy curve. All observed results for *last10runs* remain in a 20% precision range, while *oneShot* has some peaks over 50% of error. Both have a consistent tendency for under estimation.

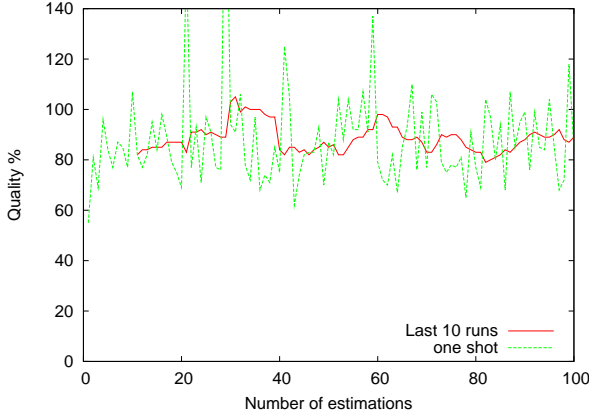


Fig. 3. HopsSampling: oneShot and last10runs heuristics, 100,000 node network

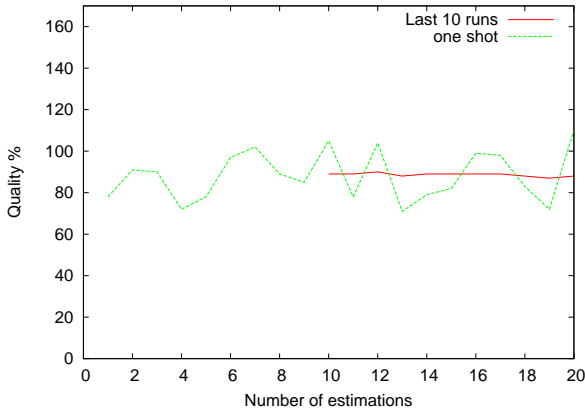


Fig. 4. HopsSampling: oneShot and last10runs heuristics, 1,000,000 node network

f) *Aggregation*: Finally, we implemented the epidemic-based aggregation protocol [9] and conducted the same set of experiments. The results match the ones presented in the original paper. One initiator node was chosen at random and initiated a gossip-based algorithm by exchanging information with a neighbor picked up at random and averaging the new value.

The results are reported on Figures 5 and 6 for respectively a 100,000 and 1,000,000 node network and give the quality of the estimation given the number of rounds. Experiments show that the size estimation naturally converges towards 100% precision around 40 rounds for 100,000 nodes and around 50 for 1,000,000 node network. The

results confirm the accuracy of the algorithm, the quality and scalability of the convergence speed in static settings.

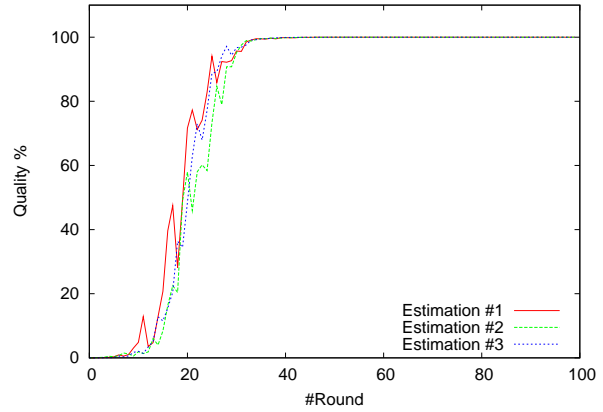


Fig. 5. Aggregation: 100,000 node network

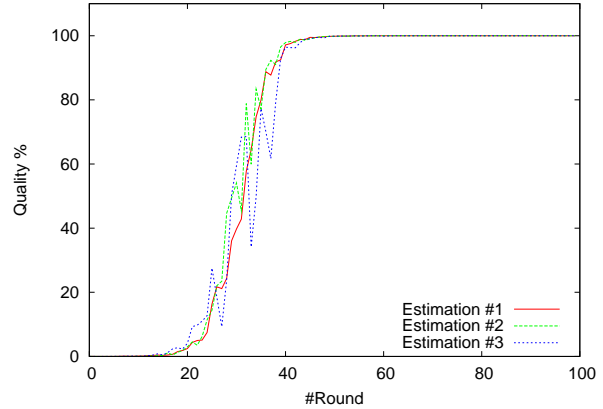


Fig. 6. Aggregation: 1,000,000 node network

g) *Scale free topology*: We now apply the counting algorithms to a 100,000 nodes scale free graph, as many large networks (like Internet) have such degree distribution properties [1], and as it seems interesting to confront these algorithms to such heterogeneous nodes' connectivity. The Figure 7 (logscale on the x and y axis) shows the power-law degree distribution of the generated graph.

The three algorithms are plotted on Figure 8; the same parameters are used: *Sample&Collide* with $l=200$, each *Aggregation* estimation occurs after 50 rounds and *HopsSampling* is used with *last10runs* heuristic.

As expected with the formal analysis in [15], the degree distribution does not bias *Sample&Collide*'s estimations. *Aggregation* also still provides accurate results. In the *HopsSampling* case, we can observe here that the under estimation factor observed earlier is increased.

h) *Static setting, summary*: Our primary conclusion is that the *Aggregation* algorithm provides accurate, and eventually exact, results in a reasonable number of cycles. This algorithm may leverage the gossip-based underlying protocol if any. In this paper, we do not take into account

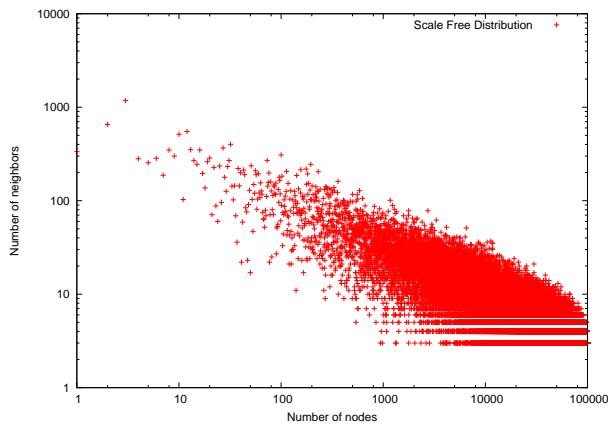


Fig. 7. Scale free degree distribution for 100,000 nodes, 3 neighbors min per node, max node degree: 1177, average: 6

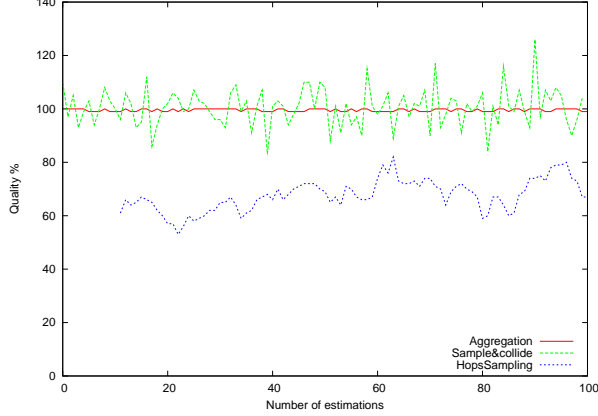


Fig. 8. Test of the 3 algorithms on a scale free graph

this underlying protocol. *Aggregation* is a sound approach for static settings, although no formal analysis is provided along with the algorithm.

HopsSampling and *Sample&Collide* produce more noisy curves. The accuracy window of *Sample&Collide* remains reasonable. We observed in the first experiments a tendency for *HopsSampling* to produce little underestimated results; this matches the trace-based simulation in [17]. The under estimate factor is amplified on the scale free graph; we address this surprising point in Section V.

D. Simulation results in dynamic context

Peer to peer overlay systems are subject to dramatic size variations over time given the versatility of the participating entities. A distributed estimation size algorithm should be able to react to such dramatic changes.

In order to simulate a changing environment, we applied constant nodes arrivals and departures (+/-50%) as well as catastrophic failures (-25%) on the simple heterogeneous random graph, and evaluated the accuracy of the size estimation in such contexts.

As the number of nodes in the system now keeps varying, the value on the y-axis of the figures is no longer

normalized but represents the actual network size.

i) *Sample&Collide*: The algorithm has to be executed perpetually in order to track size variations; the monitoring process should sample continuously the system in order to provide periodical estimations. Figures 9, 10 and 11 report the evolution of the estimation in a 100,000 node network in three scenarii: a catastrophic failure scenario, a growing scenario and a shrinking one.

All results show that the algorithm reacts very well to changes, even brutal, of sizes in the network and is therefore able to handle effectively network dynamics.

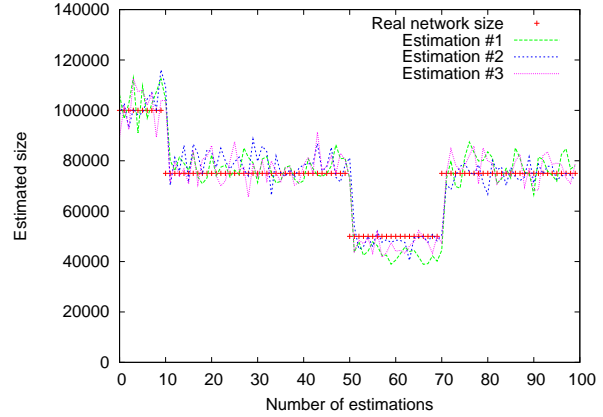


Fig. 9. *Sample&Collide*: oneShot heuristic, 100,000 node network, catastrophic failures

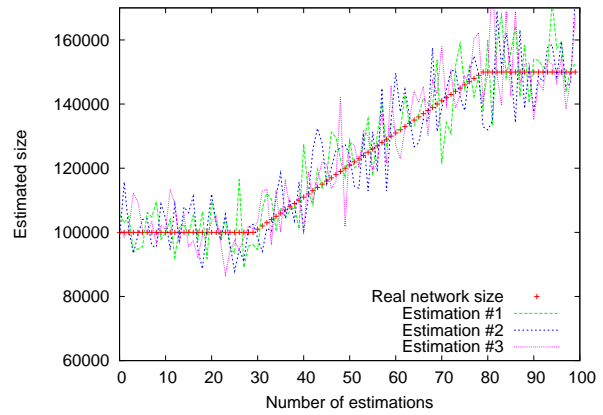


Fig. 10. *Sample&Collide*: oneShot, 100,000 node network, growing network

j) *HopsSampling*: To cope with changing scenario, the algorithm needs to be restarted periodically, there is no gradual adaptation of the estimation as new polls must be gossiped for each estimation.

Figures 12, 14 and 13 report the evolution of the estimation in a 100,000 node network, implementing the *last10runs* heuristic in the three considered scenarii. The results show a good algorithm behavior facing dynamic network, without any impact due to the overlay degradation or increase. However, we can observe that results remain

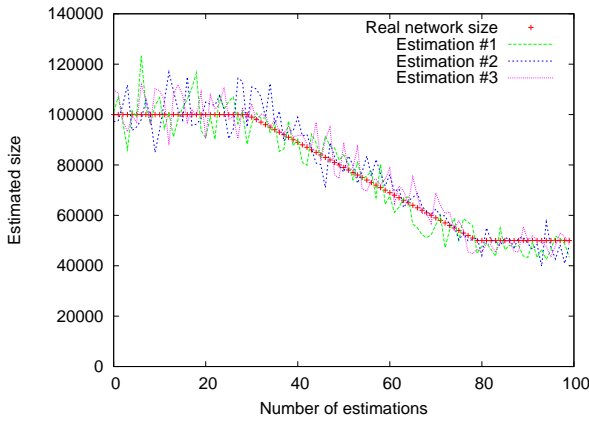


Fig. 11. *Sample&Collide*: oneShot, 100,000 node network, shrinking network

slightly under estimated and show a higher variation around the real size than in *Sample&Collide*.

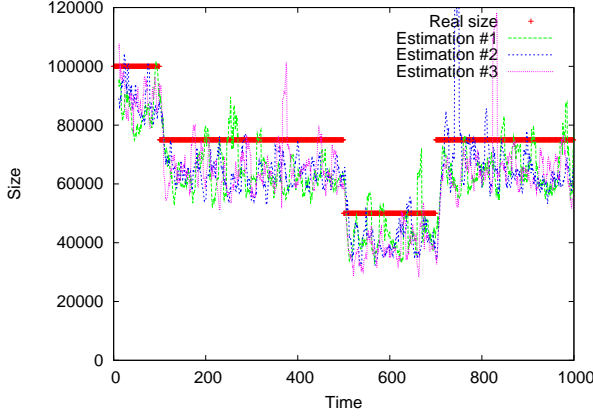


Fig. 12. *HopsSampling*: Last10runs heuristic, 100,000 node network, catastrophic failures

k) *Aggregation*: In a dynamic context, the aggregation algorithm is unable, and this is intrinsic to the way the algorithm works due to the value setting, to take node's arrivals/departures into account once an evaluation process is launched. This is due to the fact that during the aggregation process, a node sends its information to a random *alive* neighbor. Thus, there is a conservative effect, as removed nodes no longer participate and as new nodes do not get synchronized information. The overlay size estimated is only accurate at the time where the process was started.

To track size variations, the solution is to reinitialize an aggregation process at regular time intervals. By using tags (unique identifiers) on each new counting process, the algorithm can be reinitialized on demand: a node which is reached by a counting message with a new tag will create a 0 initial value and will start to participate to the active process. In our simulations, we set this time interval as the number of rounds needed by the gossip process to provide an accurate result. We observed on 5 and 6 that this value

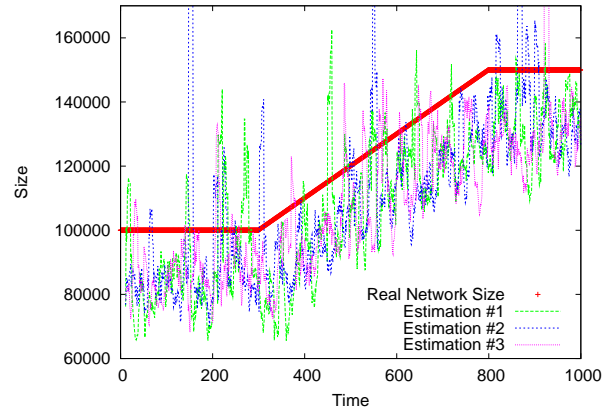


Fig. 13. *HopsSampling*: Last10runs heuristic, 100,000 node network, growing network

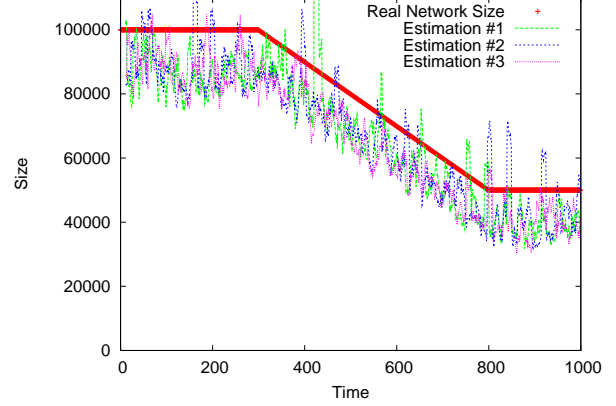


Fig. 14. *HopsSampling*: Last10runs heuristic, 100,000 node network, shrinking network

is approximately 40 for 100,000 nodes overlay, and 50 for 1,000,000, for a 99% convergence. In order not to make any hypothesis on the targeted system size, we took 50 in the following simulations (not more for a fair comparison, because this parameter is overhead consuming as shown in IV-E). This value represents the best possible algorithm's reactivity, in terms of latency, for an accurate estimation.

Figure 15, 16 and 17 report the estimation results in the three considered scenario. We observe that the algorithm does not cope well with the decrease of the network size as observed in the shrinking scenario and the catastrophic scenario. Actually, Figures 17 and 15 show a reasonable algorithm behavior until a certain threshold of nodes departures (around 30% for 100,000 nodes, with 50 rounds per estimation) is reached. We believe that this is due to the loss of connectivity of the overlay. This prevents the propagation of the estimation information across the network. The solution to correct this point is to let a larger number of gossip rounds to elapse before each estimation; this will permit the information to spread better despite low graph connectivity, but will introduce more overhead.

However, we observe that the algorithm provides a fairly good adaptation to a growing network.

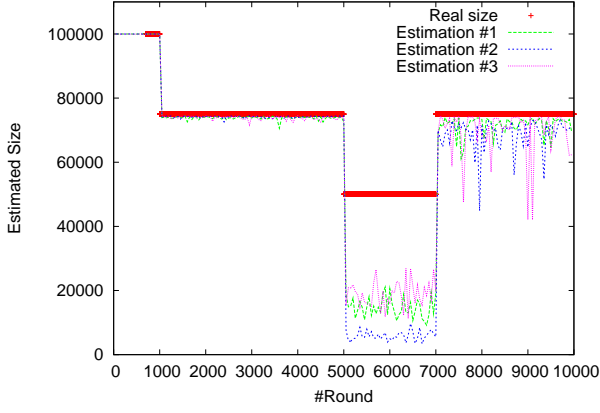


Fig. 15. Aggregation: Reaction under failures, 100,000 nodes at beginning, -25% of nodes at 100 and 500, +25,000 nodes at 700

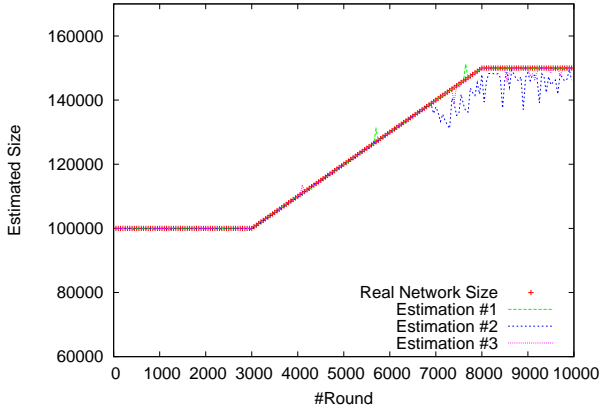


Fig. 16. Aggregation: Growing network, 100,000 node network

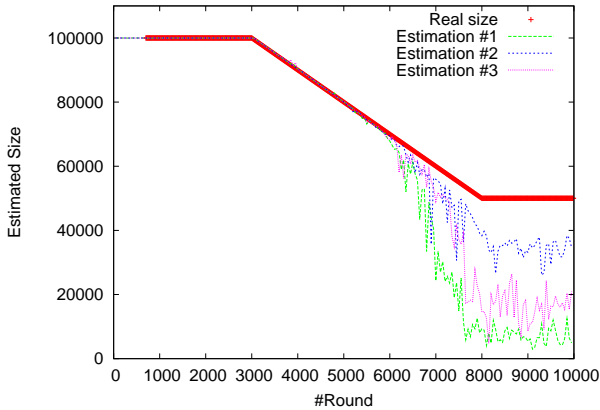


Fig. 17. Aggregation: Shrinking network, 100,000 node network

l) *Dynamic settings, summary:* We observe that *Sample&Collide* and *HopsSampling* provide fairly good results when facing brutal as well as gradual changes in the size of the system. However, *HopsSampling* as well as *Aggregation* would require an estimation process to be launched periodically.

Sample&Collide provides really reactive results; this could be explained by the *oneShot* heuristic as the algorithm does not keep any memory, while for *HopsSampling* with its *last10runs* average, there is a little convergence time to elapse to produce an accurate estimation facing a brutal topology changes.

Given that some parameters are well tuned, each of those algorithms would be able to provide accurate results. this is then a matter of acceptable overhead as discussed in the next section.

E. Algorithms overheads

Besides the accuracy of the estimation, an important parameter is the overhead induced by each method. We measure the overhead of the different algorithms as the total number of messages sent to produce the estimation. This includes spreading messages for *Aggregation* and for *HopsSampling*, return messages for *HopsSampling*, the message associated to the random walk for *Sample&Collide* as well as each sampled node's return.

- *Sample&Collide*: Analysis of costs incurred in the same simulator are provided in [15]. More specifically, we observe good scalability properties as the algorithm with $l = 100$ incurs a cost which is only 3.27 times the one incurred for $l = 10$. In order to have a more accurate result, and to compute the associated overhead, we ran some experiments with $l = 200$. The measured overhead of an estimation is around 480,000 messages, which represents 1.40 times the one incurred for $l = 100$.
- *HopsSampling*: With current parameters, a single shot estimation consumes $O(2N)$; *last10runs* thus uses $O(10 * 2N)$ messages for an estimation.
- *Aggregation*: The overhead of an estimation for *Aggregation* algorithm is fairly simple to compute: $Overhead = number_of_nodes * number_of_rounds * 2$. 2 is due to the push/pull heuristic¹.

The overhead associated to each algorithm is fairly simple to compute based on the algorithm characteristics. Table I summarizes these results and gives a overhead comparison between the three algorithms, considering the average precision of the estimation. We observe that

¹the push/pull heuristic means that when two nodes communicates, they both receive the information of the other and re-compute the current estimation

Aggregation provides a nearly perfect estimation with a high overhead compared to the other algorithms. This is the good candidate for really stringent application needs. Note that *Sample&Collide* results, when averaged with the *last10runs* heuristic and with $l = 200$, are at only few percents of the exact system size for a half overhead compared to *Aggregation*. *Sample&Collide* outperforms *HopsSampling* both in terms of the quality of the estimations and the overhead.

One advantage of *Sample&Collide* and *HopsSampling* is that they provide an opportunity to limit the overhead at the price of a less accurate estimation. On the contrary, the *Aggregation* method does not provide this flexibility and trying to obtain a rough idea of the system while keeping the overhead low could lead to really bad results.

V. Comparisons and tradeoffs

m) Tradeoffs: With respect to tradeoffs, we observed that *Sample&Collide* provides the most flexible algorithm and offers a wide range of possible configurations depending on the accuracy of the estimate and associated overhead. For example, a really accurate result might not always be important and could furthermore uselessly waste resources. In this case, *Sample&Collide* with a smaller l parameter ($l = 10$ for example) could be a cheap solution for a relatively good result, as depicted on Figure 18 where on average only 100,000 messages are used for an estimation. Our experiments show that for *HopsSampling*, using a lower *minHopsReporting* parameter does not significantly reduce the overhead, while degrading accuracy. *HopsSampling*, to a lower extent also provides some flexibility. However *Aggregation* does not, the only option is a high precision associated to a fairly large overhead. As detailed in [15], *Sample&Collide* could as well produced really precise results by increasing l parameter, and hence could compete with aggregation. A strength of this algorithm is thus to adapt to the application performance needs by simply modifying one parameter.

n) Dynamicity: With respect to dynamic networks, we observed that *Sample&Collide* provides reliable results despite the degradation of the overlay connectivity. On the other hand, gossip-based algorithms need more time to complete a run due to the increase of the epidemic propagation time. The *HopsSampling* algorithm does not suffer from the same problem as it sends more messages each round, hence increasing the propagation efficiency. *Aggregation* and *HopsSampling* imply that each node in the overlay participates to the current counting process, which systematically entails resources consumption on each node at each process round.

o) Accuracy: In terms of accuracy, *Aggregation* outperforms the other algorithms and provides in a reasonable

number of rounds an almost exact estimate. This approach should be used for applications with the most stringent needs.

The polling algorithm of *HopsSampling* is not biased: we verified our intuition by giving the accurate distance from the initiator to all nodes in the overlay, and the resulting size estimation was correct. The under estimation phenomenon thus concerns the first phase of the algorithm (the distance computation) certainly due to the fact that not all nodes are reached during the process message spread (approximately 11% of non reached nodes out of 100,000), and because the distances from the node initiator are not always accurate. At the time we are writing this paper, we are looking to see if this message spread problem is due to non optimal parameter setting, or algorithm/simulator imprecision. Basic broadcast techniques, at the price of a higher overhead, are probably an easy solution to this issue.

p) Other points: *HopsSampling* has the drawback of creating a message flood towards the initiator during the *HopCount* collection process; this produces a sensitive point in the network, and may overload the initiator's neighbors. *Sample&Collide*, which provides less overhead and more accurate results, might be preferred by the user. However, *HopsSampling* probably outperforms the other algorithms in terms of delay, which we haven't measured in this comparison due to the fact that physical network topology was not modeled in our simulator. A gossip based broadcast and an immediate ACK response from the nodes that succeeds the probabilistic test is very likely to be much shorter than the 50 rounds of *Aggregation* or the wait for 200 equivalent samples of *Sample&Collide*.

An interesting characteristic of the *Aggregation* is that eventually the size estimation is available at each node of the network, as every overlay node possesses the local information to compute the estimation. There is no need for a broadcast of the size estimation as in the two other algorithms where only the initiator gets to compute the size estimation.

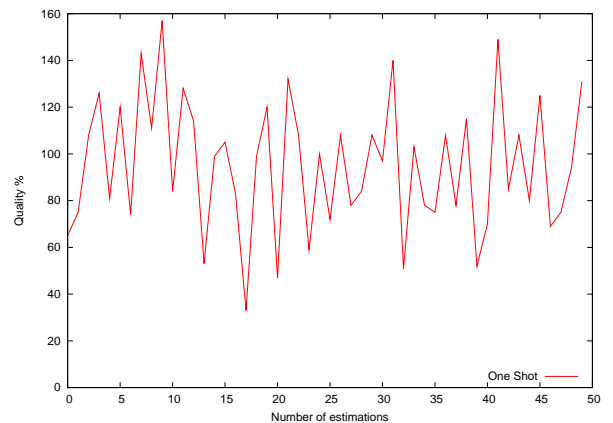


Fig. 18. Sample & collide with $l=10$, 100,000 node network

Algorithm	Sample&Collide (l=200)	HopsSampling	Sample&Collide (l=200)	Aggregation
Parameters	oneShot	last10runs	last10runs	50 rounds
Accuracy	+/- 10%	- 20%	+/- 4%	- 1%
Overhead	0,5M	2,5M	5M	10M

TABLE I. Example of algorithm's overhead for an estimation (in millions of messages) on a 100,000 node overlay

VI. Conclusion

Monitoring a peer to peer system is a core functionality of distributed system and might be used in many applicative contexts. More specifically, the system size is very often needed to set the parameter values or monitor a system behavior and this can be used for a wide spectrum of applications ranging from sensor networks to Grid computing.

In this paper, we provide the first comparison of some solutions computing the size of a large-scale system in a fully decentralized manner, in peer to peer unstructured overlays where peers have only a restricted knowledge of the system. We evaluated three different algorithms, representative of the probabilistic polling, random walk-based and gossip-based approaches and compared them along their accuracy, convergence, capability to handle network size changes and the associated overhead.

We observe that all three algorithms have passed the scalability and dynamicity tests but significant differences arised. Not surprisingly, as in many areas of computer science, there is a natural tradeoff between the quality of the estimate and the associated overhead, here bandwidth and computational resources. With respect to this tradeoff, *Sample&Collide* seems to be the most flexible algorithm.

As part of future work, the physical network modeling would be an interesting goal and might provide new insights on the comparison.

References

- [1] R. Albert and A.-L. Barabási. Statistical mechanics of complex networks. *Reviews of Modern Physics* 74, 47 (2002), 2002.
- [2] M. Bawa, H. Garcia-Molina, A. Gionis and R. Motwani. Estimating aggregates on a peer-to-peer network. *Technical Report, Dept. of computer science, Stanford University*, 2003.
- [3] M. Castro, P. Druschel, A. Ganesh, A. Rowstron, and D. S. Wallach. Security for structured peer-to-peer overlay networks. *OSDI 02*, December 2002.
- [4] P. Eugster, S. Handurukande, R. Guerraoui, A.-M. Kermarrec and P. Kouznetsov. Lightweight probabilistic broadcast. *ACM Transaction on Computer Systems*, 21(4), November 2003.
- [5] M. Frans Kaashoek and David R. Karger. Koorde: A simple degree-optimal distributed hash table. *IPTPS 2003*, February 2003.
- [6] T. Friedman and D. Towsley. Multicast session membership size estimation. *IEEE INFOCOMM '99*, 1999.
- [7] A. J. Ganesh, A.M. Kermarrec and L. Massoulié Peer-to-peer membership management for gossip-based protocols. *IEEE transactions on computers*, vol.52, no.2, February 2003, 2003.

- [8] M. Jelasity, R. Guerraoui, A.-M. Kermarrec, and M. van Steen. The Peer Sampling Service: Experimental Evaluation of Unstructured Gossip-Based Implementations. *Middleware 2004*, October 2004.
- [9] M. Jelasity and A. Montresor. Epidemic-style proactive aggregation in large overlay networks. *ICDCS 2004*, 2004.
- [10] M. Jelasity, S. Voulgaris, R. Guerraoui, A.-M. Kermarrec. Gossip-based peer sampling. *Submitted*.
- [11] D. Kostoulas, D. Psaltoulis, I. Gupta, K. Birman, and A. Demers. Decentralized schemes for size estimation in large and dynamic groups. *IEEE NCA '05*, 2005.
- [12] D. Malkhi, M. Naor and D. Ratajczak. Viceroy: a scalable and dynamic emulation of the butterfly. *PODC 2002*, July 2002.
- [13] G. S. Manku. Routing networks for distributed hash tables. *PODC 2003*, June 2003.
- [14] G. S. Manku, M. Bawa and P. Raghavan, Symphony: Distributed Hashing in a Small World *USITS 2003*, 2003.
- [15] L. Massoulié, E. Le Merrer, A.-M. Kermarrec and A.J. Ganesh. Peer counting and sampling in overlay networks: random walk methods. *To appear in PODC 2006*, July 2006.
- [16] Dimitrios Psaltoulis. *Private communication.*, July 2005.
- [17] D. Psaltoulis, D. Kostoulas, I. Gupta, K. Birman, and A. Demers. Practical algorithms for size estimation in large and dynamic groups. *PODC 2004*, 2004.
- [18] R. van Renesse, K. Birman and W. Vogels. Astrolabe: A robust and scalable technology for distributed system monitoring, management, and data mining. *ACM Transactions on Computer Systems (TOCS)*, Vol. 21 pp. 164-206, May 2003.
- [19] S. Voulgaris, D. Gavidia, M. van Steen. CYCLON: Inexpensive Membership Management for Unstructured P2P Overlays. *Journal of Network and Systems Management*, Vol. 13, No. 2, June 2005.